

What Is Claimed Is:

1. A method for communicating between applications executing in different partitions of a partitionable computer system, the method comprising:

5 receiving a request made by a first application on a first partition to establish a network connection with a second application on a second partition and to send a message to the second application via the network connection;

10 establishing a connection between the first partition and the second partition of the computer system through a memory region of the computer system shared by both the first partition and the second partition, wherein the connection emulates the requested network connection; and

15 sending the message to the second application via the connection established through the shared memory region, whereby the connection established through the shared memory region appears to the first and second applications as the requested network connection.

20 2. The method recited in claim 1, wherein the connection requested by the first application comprises a socket connection, and wherein the step of establishing a connection through the shared memory region comprises establishing a connection through the shared memory region that emulates a socket connection.

25 3. The method recited in claim 1, wherein the steps of establishing and sending comprise:

creating a data structure in the shared memory region comprising a plurality of data segments forming a circular buffer;

writing, from the first partition on behalf of the first application, the message to one or more data segments, as needed, and updating an indication of the data segment containing the most recently written portion of the message.

30 4. The method of claim 3, wherein updating an indication of the data segment containing the most recently written portion of the message comprises incrementing a head index.

5. The method recited in claim 3, further comprising:

reading, from the second partition on behalf of the second application, the message from said one or more data segments and updating an indication of which data segments have been read from the data structure; and

5 providing the message read from the data structure to the second application in accordance with an API associated with the requested network connection.

6. The method of claim 5, wherein updating an indication of which data segments have been read from the data structure comprises incrementing a tail index.

10 7. The method of claim 1, further comprising polling, by the receiving partition, the shared memory region to determine if the message has been written to the shared memory region.

15 8. The method of claim 1, further comprising receiving, by the receiving partition, an interrupt initiated by the sending partition and indicating that the message has been written to the shared memory region.

20 9. The method of claim 2, wherein the step of establishing a connection between the first partition and the second partition of the computer system that emulates a socket connection further comprises performing the following steps on the second partition:

(a) creating a socket on behalf of the second application to listen for attempts to connect thereto;

(b) receiving a connect message from the first partition that identifies a memory location of the shared memory region at which the first partition has allocated a first data area to serve as a buffer for transferring data from the first partition to the second partition;

(c) matching the received connect message to the listening socket created in step (a);

(d) allocating a second data area in the shared memory region to serve as a buffer for transferring data from the second partition to the first partition;

30 (e) mapping both the first and second data areas into a process space of the listening socket;

(f) initializing the second data area; and

(g) returning a "connected" indication to the first partition and informing the application on the second partition that the socket connection has been established.

10. The method of claim 9, further comprising performing the following steps
5 on the first partition:

(a') receiving the request from the first application to establish the socket connection with the second application;

(b') creating a connecting socket;

(c') allocating the first data area in the shared memory region;

10 (d') sending the connect message to the second partition that identifies the memory location of the shared memory region at which the first data area has been allocated; and

(e') upon receipt of the "connected" indication from the second partition, mapping the first and second data areas into a process space of the connecting socket to establish the socket connection between the first and second partitions.

15 11. The method of claim 1, further comprising:

creating, in the shared memory region, a plurality of output queues, one for each of said first and second partitions, the output queue for a given partition indicating whether that partition has placed in the shared memory window a message intended for any of the other partitions and if so, identifying a buffer containing the message, each partition polling the output queues of the other partitions to determine whether those other partitions have placed any messages intended for it in the shared memory window;

receiving at the first partition from the first application, said request to send a message to the second application via the requested type of network connection;

25 writing, in response to the received request, the message to an available buffer in the shared memory window and indicating in the output queue of the first partition that the message has been written thereto;

determining, at the second partition, from the output queue of the first partition, that the message has been placed in the buffer and retrieving the message from the buffer; and

30 providing the message read from the data structure to the second application in accordance with an API associated with the requested network connection.

12. A computer readable medium having program code store thereon for communicating between applications executing in different partitions of a partitionable computer system, the program code, when executed on a processor, causing the processor to perform the following:

5 receiving a request made by a first application on a first partition to establish a network connection with a second application on a second partition and to send a message to the second application via the network connection;

10 establishing a connection between the first partition the second partition of the computer system through a memory region of the computer system shared by both the first partition and the second partition, wherein the connection emulates the requested network connection; and

15 sending the message to the second application via the connection established through the shared memory region, whereby the connection established through the shared memory region appears to the first and second applications as the requested network connection.

20 13. The computer readable medium recited in claim 12, wherein the connection requested by the first application comprises a socket connection, and wherein the step of establishing a connection through the shared memory region comprises establishing a connection through the shared memory region that emulates a socket connection.

25 14. The computer readable medium recited in claim 12, wherein the steps of establishing and sending comprise:

creating a data structure in the shared memory region comprising a plurality of data segments forming a circular buffer;

writing, from the first partition on behalf of the first application, the message to one or more data segments, as needed, and updating an indication of the data segment containing the most recently written portion of the message.

30 15. The computer readable medium of claim 14, wherein updating an indication of the data segment containing the most recently written portion of the message comprises incrementing a head index.

16. The computer readable medium recited in claim 14, wherein the program code, when executed on a processor, further causes the processor to perform the following:

reading, from the second partition on behalf of the second application, the message from said one or more data segments and updating an indication of which data segments have been read from the data structure; and

providing the message read from the data structure to the second application in accordance with an API associated with the requested network connection.

17. The computer readable medium of claim 16, wherein updating an indication of which data segments have been read from the data structure comprises incrementing a tail index.

18. The computer readable medium of claim 12, wherein the program code, when executed on a processor, further causes the processor to poll, by the receiving partition, the shared memory region to determine if the message has been written to the shared memory region.

19. The computer readable medium of claim 12, wherein the program code, when executed on a processor, further causes the processor to receive, by the receiving partition, an interrupt initiated by the sending partition and indicating that the message has been written to the shared memory region.

20. The method of claim 13, wherein the step of establishing a connection between the first partition and the second partition of the computer system that emulates a socket connection further comprises performing the following steps on the second partition:

25 (a) creating a socket on behalf of the second application to listen for attempts to connect thereto;

(b) receiving a connect message from the first partition that identifies a memory location of the shared memory region at which the first partition has allocated a first data area to serve as a buffer for transferring data from the first partition to the second partition;

30 (c) matching the received connect message to the listening socket created in step (a);

(d) allocating a second data area in the shared memory region to serve as a

buffer for transferring data from the second partition to the first partition;

- (e) mapping both the first and second data areas into a process space of the listening socket;
- (f) initializing the second data area; and
- 5 (g) returning a "connected" indication to the first partition and informing the application on the second partition that the socket connection has been established.

21. The method of claim 20, wherein the step of establishing a connection further comprises performing the following steps on the first partition:

- 10 (a') receiving the request from the first application to establish the socket connection with the second application;
- (b') creating a connecting socket;
- (c') allocating the first data area in the shared memory region;
- (d') sending the connect message to the second partition that identifies the memory location of the shared memory region at which the first data area has been allocated; and
- 15 (e') upon receipt of the "connected" indication from the second partition, mapping the first and second data areas into a process space of the connecting socket to establish the socket connection between the first and second partitions.

20 22. The computer readable medium of claim 12, wherein the program code, when executed on a processor, further causes the processor to perform the following:

25 creating, in the shared memory region, a plurality of output queues, one for each of said first and second partitions, the output queue for a given partition indicating whether that partition has placed in the shared memory window a message intended for any of the other partitions and if so, identifying a buffer containing the message, each partition polling the output queues of the other partitions to determine whether those other partitions have placed any messages intended for it in the shared memory window;

receiving at the first partition from the first application, said request to send a message to the second application via the requested type of network connection;

30 writing, in response to the received request, the message to an available buffer in the shared memory window and indicating in the output queue of the first partition that the message has been written thereto;

determining, at the second partition, from the output queue of the first partition, that the message has been placed in the buffer and retrieving the message from the buffer; and providing the message read from the data structure to the second application in accordance with an API associated with the requested network connection.

5

23. A computer system comprising:

a plurality of processing modules, groups of one or more processing modules being configured as separate partitions within the computer system, each partition operating under the control of a separate operating system;

10 a main memory to which each processing module is connected, the main memory having defined therein at least one shared memory region to which at least two different ones of said separate partitions have shared access; and

15 program code, executing on each of at least a first partition and a second partition of the computer system, which program code establishes a connection between a first application on the first partition of the computer system and a second application on the second partition of the computer system through the shared memory region, wherein the connection through the shared memory region emulates a network connection requested by one of the applications.

20 24. The computer system recited in claim 23, wherein said program code executing on each of said first and second partitions comprises a shared memory service provider that serves as an interface between a component of the computer system that provides an API through which said application can make said request for a network connection and the shared memory region of the main memory through which the emulated network connection is established.

25

25. The computer system recited in claim 24, wherein the shared memory service provider on each of said first and second partitions establishes a data structure in the shared memory region through which data is transferred from that partition to the shared memory service provider on the other partition.

30

26. The computer system recited in claim 25, wherein the data structure comprises:

a plurality of data segments, each of the plurality of data segments for storing network message data to be sent from a sending shared memory service provider to a receiving shared memory service provider;

5 a control segment for controlling reading and writing of data in the plurality of data segments, the control segment comprising:

 a first portion comprising:

 a first field for storing an indication of the data segment containing the most recently written network message data; and

10 a second field for storing an indication of the data segment containing the earliest written, but not read, network message data;

 and

 a plurality of second portions, each second portion corresponding to one of the plurality of data segments for control of the data segment, each second portion comprising:

15 a first field for storing an indication of the beginning of network message data within the data segment; and

 a second field for storing an indication of the end of network message data within the data segment.

20 27. The computer system recited in claim 26, wherein the first portion further comprises:

 a third field for storing an indication that the sending shared memory service provider is waiting to send the network message; and

25 a fourth field for storing an indication that the receiving shared memory service provider is waiting to receive the network message.

28. The computer system recited in claim 26, wherein each second portion further comprises a third field for storing an indication of a length of network message data within the data segment.

30

29. The computer system recited in claim 28, wherein the plurality of data segments are linked to form a circular buffer, and wherein each second portion further comprises:

a fourth field for storing an indication of the next data segment in the circular buffer; and

a fifth field for storing an indication that the data segments contains a last portion of a network message stored across a plurality of data segments.

5

30. The computer system recited in claim 25, wherein the computer system provides a resource through which the shared memory service provider can establish the data structure and control the transfer of data through it, the resource providing the ability to perform at least one of the following operations on the shared memory region: (i) allocate an area of the shared memory region; (ii) map and unmap an area of the shared memory region; deallocate an area of the shared memory region; (iii) send and receive signals to and from other partitions via the shared memory region; and (iv) receive status information about the shared memory region and about selected partitions.

15 31. The computer system recited in claim 24, wherein the shared memory service provider comprises:

20 a dynamic link library (DLL) that executes in a user mode of the operating system of its respective partition, there being an instance of the shared memory service provider DLL in a process space of each application in the partition that may request the establishment of a network connection; and

a device driver that executes in a kernel mode of the operating system of the respective partition, there being only one instance of the device driver in each partition.

25 32. The computer system recited in claim 23, wherein the connection established through the shared memory region emulates a socket connection.

33. The computer system recited in claim 32, wherein said program code executing on each of said at least first and second partitions comprises a shared memory service provider that serves as an interface between a component of the computer system that provides an API through which an application can make a request for a socket connection and the shared memory region of the main memory through which the emulated socket connection is established.

34. The computer system recited in claim 33, wherein the operating system in each partition comprises a MICROSOFT WINDOWS operating system, and wherein the component of the computer system that provides the API of the requested socket connection comprises a Winsock DLL and a Winsock Switch, the Winsock DLL forwarding a request for a socket connection made by an application in a given partition to the Winsock Switch, which Winsock Switch allows multiple service providers, each of which provide TCP/IP services, to service such a request, and wherein the shared memory service provider acts as a TCP/IP service provider so that a request from an application for a socket connection can be serviced by the shared memory service provider.

35. The computer system recited in claim 34, wherein the shared memory service provider on a first partition that represents the listening side of a requested socket connection performs the following steps:

- (a) creating a socket on behalf of a first application executing in the first partition in order to listen for attempts to connect thereto;
- (b) receiving a connect message from the shared memory service provider on the second partition that identifies a memory location of the shared memory region at which the shared memory service provider on the second partition has allocated a first data area to serve as a buffer for transferring data from the second partition to the shared memory service provider on the first partition;
- (c) matching the received connect message to the listening socket created in step (a);
- (d) allocating a second data area in the shared memory region to serve as a buffer for transferring data from the first partition to the shared memory service provider on the second partition;
- (e) mapping both the first and second data areas into a process space of the listening socket;
- (f) initializing the second data area; and
- (g) returning a "connected" indication to the shared memory service provider on the second partition and informing the application on the first partition that a socket connection has been established.

36. The computer system recited in claim 35, wherein the shared memory service provider on the second partition performs the following steps:

- (a') receiving a request from an application on the second partition to establish a socket connection with the first application on the first partition;
- (b') creating a connecting socket on the second partition;
- (c') allocating the first data area in the shared memory region;
- (d') sending the connect message to the first partition that identifies the memory location of the shared memory region at which the first data area has been allocated; and
- (e') upon receipt of the "connected" indication from the first partition, mapping the first and second data areas into a process space of the connecting socket to establish the socket connection between the first and second partitions.

37. The computer system recited in claim 23, wherein the program code implements a polling process by which each partition polls an area within the shared memory region to determine whether any communications intended for it have been placed in the shared memory region by another partition.

38. The computer system recited in claim 37, wherein the area comprises a plurality of output queues, one for each partition, the output queue for a given partition indicating whether that partition has placed in the shared memory region any communications intended for any of the other partitions, each partition polling the output queues of the other partitions to determine whether those other partitions have placed any communications intended for it in the shared memory region.

39. The computer system recited in claim 38, wherein for any communications placed in the shared memory region by a sending partition and intended to be received by another partition, the output queue of the sending partition specifies the location within the shared memory region of a buffer containing that communication.

40. The computer system recited in claim 39, wherein the program code executing on each of said first and second partitions further comprises a shared memory driver

that receives a request to send a message to an application on another partition, the request having been made in accordance with the application programming interface (API) associated with the requested type of network connection, and that, in response to the request, causes the message to be placed in an available buffer in the shared memory region and causes an indication 5 of the message to be placed in the output queue of the sending partition.

41. The computer system recited in claim 40, wherein the shared memory driver on each partition implements a same interface as a network device driver to enable application programs and the operating system on that partition to send communications to other 10 partitions via the shared memory region in the same manner that communications are sent to other computer systems over a network via a network interface card.

42. A computer-readable medium having stored thereon a data structure for emulating network communications between a sending program and a receiving program executing on a computer system, the data structure comprising:

a plurality of data segments, each of the plurality of data segments for storing a network message;

a control segment for controlling reading and writing of data in the plurality of data segments, the control segment comprising:

a first portion including:

a first field for storing an indication of the data segment containing the most recently written network message; and

a second field for storing an indication of the data segment containing the earliest written, but not read, network message;

25 and

a plurality of second portions, each second portion corresponding to one of the plurality of data segments for control of the data segment, each second portion comprising:

30 a first field for storing an indication of the beginning of network message data within the data segment; and

a second field for storing an indication of the end of network message data within the data segment.

43. The computer-readable medium of claim 42, wherein the first portion further comprises:

- 5 a third field for storing an indication that the sending program is waiting to send the network message; and
- a fourth field for storing an indication that the sending program is waiting to receive the network message.

44. The computer-readable medium of claim 42, wherein each second portion 10 further comprises:

- a third field for storing an indication of a length of network message data within the data segment.

45. The computer-readable medium of claim 44, wherein the plurality of data segments are linked to form a circular buffer, and wherein each second portion further comprises:

15 a fourth field for storing an indication of the next data segment in the circular buffer; and

 a fifth field for storing an indication that the data segment contains a last portion of a network message stored across a plurality of data segments.

20 46. Apparatus of use in a partitionable computer system that comprises a plurality of processing modules, groups of one or more processing modules being configured as separate partitions within the computer system, each partition operating under the control of a separate operating system, and wherein the computer system further comprises a main memory to which each processing module is connected, the main memory having defined therein at least one shared memory region to which at least two different ones of said separate partitions have shared access, said apparatus comprising:

- 25 program code, executing on each of at least a first partition and a second partition of the computer system, which program code establishes a connection between a first application on the first partition of the computer system and a second application on the second partition of the computer system through the shared memory region, wherein the connection through the shared memory region emulates a network connection requested by one of the applications.

47. The apparatus recited in claim 46, wherein said program code executing on each of said first and second partitions comprises a shared memory service provider that serves as an interface between a component of the computer system that provides an API through which 5 said application can make said request for a network connection and the shared memory region of the main memory through which the emulated network connection is established.

48. The apparatus recited in claim 47, wherein the shared memory service provider on each of said first and second partitions establishes a data structure in the shared 10 memory region through which data is transferred from that partition to the shared memory service provider on the other partition.

49. The apparatus recited in claim 48, wherein the data structure comprises:
15 a plurality of data segments, each of the plurality of data segments for storing network message data to be sent from a sending shared memory service provider to a receiving shared memory service provider;
a control segment for controlling reading and writing of data in the plurality of data segments, the control segment comprising:
20 a first portion comprising:
a first field for storing an indication of the data segment containing the most recently written network message data; and
a second field for storing an indication of the data segment containing the earliest written, but not read, network message data;
and
25 a plurality of second portions, each second portion corresponding to one of the plurality of data segments for control of the data segment, each second portion comprising:
a first field for storing an indication of the beginning of network message data within the data segment; and
30 a second field for storing an indication of the end of network message data within the data segment.

50. The apparatus recited in claim 49, wherein the first portion further comprises:

a third field for storing an indication that the sending shared memory service provider is waiting to send the network message; and

5 a fourth field for storing an indication that the receiving shared memory service provider is waiting to receive the network message.

51. The apparatus recited in claim 49, wherein each second portion further comprises a third field for storing an indication of a length of network message data within the
10 data segment.

52. The apparatus recited in claim 51, wherein the plurality of data segments are linked to form a circular buffer, and wherein each second portion further comprises:

15 a fourth field for storing an indication of the next data segment in the circular buffer; and

a fifth field for storing an indication that the data segments contains a last portion of a network message stored across a plurality of data segments.

53. The apparatus recited in claim 48, wherein the computer system provides a resource through which the shared memory service provider can establish the data structure and control the transfer of data through it, the resource providing the ability to perform at least one of the following operations on the shared memory region: (i) allocate an area of the shared memory region; (ii) map and unmap an area of the shared memory region; deallocate an area of the shared memory region; (iii) send and receive signals to and from other partitions via the shared memory region; and (iv) receive status information about the shared memory region and about selected partitions.
20
25

54. The apparatus recited in claim 47, wherein the shared memory service provider comprises:

30 a dynamic link library (DLL) that executes in a user mode of the operating system of its respective partition, there being an instance of the shared memory service provider DLL in a process space of each application in the partition that may request the establishment of a

network connection; and

a device driver that executes in a kernel mode of the operating system of the respective partition, there being only one instance of the device driver in each partition.

5 55. The apparatus recited in claim 46, wherein the connection established through the shared memory region emulates a socket connection.

10 56. The apparatus recited in claim 55, wherein said program code executing on each of said at least first and second partitions comprises a shared memory service provider that serves as an interface between a component of the computer system that provides an API through which an application can make a request for a socket connection and the shared memory region of the main memory through which the emulated socket connection is established.

15 57. The apparatus recited in claim 56, wherein the operating system in each partition comprises a MICROSOFT WINDOWS operating system, and wherein the component of the computer system that provides the API of the requested socket connection comprises a Winsock DLL and a Winsock Switch, the Winsock DLL forwarding a request for a socket connection made by an application in a given partition to the Winsock Switch, which Winsock Switch allows multiple service providers, each of which provide TCP/IP services, to service such a request, and wherein the shared memory service provider acts as a TCP/IP service provider so that a request from an application for a socket connection can be serviced by the shared memory service provider.

20 58. The apparatus recited in claim 57, wherein the shared memory service provider on a first partition that represents the listening side of a requested socket connection performs the following steps:

- (a) creating a socket on behalf of a first application executing in the first partition in order to listen for attempts to connect thereto;
- (b) receiving a connect message from the shared memory service provider on the second partition that identifies a memory location of the shared memory region at which the shared memory service provider on the second partition has allocated a first data area to serve as a buffer for transferring data from the second partition to the shared memory service provider on

the first partition;

- (c) matching the received connect message to the listening socket created in step (a);
- (d) allocating a second data area in the shared memory region to serve as a buffer for transferring data from the first partition to the shared memory service provider on the second partition;
- 5 (e) mapping both the first and second data areas into a process space of the listening socket;
- (f) initializing the second data area; and
- 10 (g) returning a "connected" indication to the shared memory service provider on the second partition and informing the application on the first partition that a socket connection has been established.

15 59. The apparatus recited in claim 58, wherein the shared memory service provider on the second partition performs the following steps:

- (a') receiving a request from an application on the second partition to establish a socket connection with the first application on the first partition;
- (b') creating a connecting socket on the second partition;
- (c') allocating the first data area in the shared memory region;
- 20 (d') sending the connect message to the first partition that identifies the memory location of the shared memory region at which the first data area has been allocated; and
- (e') upon receipt of the "connected" indication from the first partition, mapping the first and second data areas into a process space of the connecting socket to establish the socket connection between the first and second partitions.

25 60. The apparatus recited in claim 46, wherein the program code implements a polling process by which each partition polls an area within the shared memory region to determine whether any communications intended for it have been placed in the shared memory region by another partition.

30 61. The apparatus recited in claim 60, wherein the area comprises a plurality of output queues, one for each partition, the output queue for a given partition indicating whether

that partition has placed in the shared memory region any communications intended for any of the other partitions, each partition polling the output queues of the other partitions to determine whether those other partitions have placed any communications intended for it in the shared memory region.

5

62. The apparatus recited in claim 61, wherein for any communications placed in the shared memory region by a sending partition and intended to be received by another partition, the output queue of the sending partition specifies the location within the shared memory region of a buffer containing that communication.

10

63. The apparatus recited in claim 62, wherein the program code executing on each of said first and second partitions further comprises a shared memory driver that receives a request to send a message to an application on another partition, the request having been made in accordance with the application programming interface (API) associated with the requested type of network connection, and that, in response to the request, causes the message to be placed in an available buffer in the shared memory region and causes an indication of the message to be placed in the output queue of the sending partition.

15

64. The apparatus recited in claim 63, wherein the shared memory driver on each partition implements a same interface as a network device driver to enable application programs and the operating system on that partition to send communications to other partitions via the shared memory region in the same manner that communications are sent to other computer systems over a network via a network interface card.

20
25